

0945000-12000

5 Generally, the present invention relates to memory system architecture and,
in particular, the present invention relates to cache design.

The speed at which computer processors can execute instructions continues to outpace the ability of computer memory systems to supply instructions and data to the processors. Consequently, many high-performance computing systems provide a high-speed buffer storage unit, commonly called a cache or cache memory, between the working store or memory of the central processing unit (“CPU”) and the main memory.

A level 1 cache (“L1”) generally refers to a memory bank built closest to the
30 central processing unit (“CPU”) chip, typically on the same chip die. A level 2

Brief Description of the Drawings

Figure 1 shows a block diagram of a system in which example embodiments of the invention can be implemented

Figures 2A-2D show flow charts of example embodiments of a method of speculative execution.

Figure 3 shows a flow chart of an example embodiment of a method of replacing cache lines during run-ahead execution.

Figures 4A-4C show flow charts of example embodiments of a method of accessing a cache.

Figure 5 shows a flow chart of an example embodiment of a method of executing a software prefetching thread on a multithreaded processor.

Figures 6A-6C show block diagrams of example embodiments of a processor.

Figures 7-10 show block diagrams of example embodiments of a multiprocessor computer system.

Detailed Description

A novel method and apparatus are described for protecting cache lines from premature eviction after the cache lines were allocated by a run-ahead prefetcher. In the following detailed description of the invention reference is made to the accompanying drawings which form a part hereof, and in which is shown, by way of illustration, specific embodiments in which the invention may be practiced. In the drawings, like numerals describe substantially similar components throughout the several views. These embodiments are described in sufficient detail to enable those skilled in the art to practice the invention. Other embodiments may be utilized and structural, logical, and electrical changes may be made without departing from the scope of the present invention.

The apparatus and method of example embodiments of the invention prevents premature eviction of cache lines allocated to a cache by a run-ahead

Table 1

```

5  struct cache_line_struct{
    unsigned long tag; /* line tag */
    char valid; /* valid bit */
    char dirty; /* dirty bit */
    char protected; /* protection bit */
10  char *data; /* line data */
    } line;

    struct cache_set_struct{
        line lines[NUM_ASSOC]; /* lines[0] = LRU, lines[NUM_ASSOC-1] = MRU */
15  } cache_set;

    struct cache_struct{
        cache_set sets[NUM_CACHE_SETS];
    } cache;
20  cache c;

    char* /* return line data */
    do_cache_access(unsigned long addr, /* address of access */
25  int TOA, /* type of access RD/WR */
        int run_ahead, /* 1 = run ahead mode, 0 = normal mode */
        char* data /* for writes */
        ) {
        unsigned long tag = GET_TAG(addr);
30  unsigned int set = GET_SET(addr);
        unsigned line *l = find_line_in_set(tag, c.sets[set]);
        unsigned line *repl;

        if (!run_ahead) {
35  if (l) { /* if a hit */
            l->protected = 0;
            update_LRU(c.sets[set], l); /* place l at the head of LRU list */
            if (TOA == RD) /* read */
                return l->data;
40  else { /* write */
                l->dirty = 1;
                return l->data = data;
            }
        } else { /* miss */
45  repl = &(c.sets[set].lines[0]); /* replace LRU block */
            process_miss(addr, TOA, run_ahead, data.repl);

```


cache controller 716 is associated with the plurality of caches 708 (shown in Figures 7 and 8). In this embodiment, the control logic 714 resides in the at least one cache controller 716. However, all or part of the control logic 714 may reside elsewhere.

In one embodiment, the multiprocessor computer system 700 further comprises a plurality of tag arrays 718, as shown in Figure 10. A tag array 718 is associated with each cache 708 (shown in Figures 7 and 8). In this embodiment, the protection bits 712 reside in each tag array 718 and are associated with cache lines 710. A tag is the remainder of an address generated by the processor after the set bits have been removed. Set bits are the address used to find a line within a cache. The cache management logic may compare the tag bits of the address with the tag bits of the cache directory which are stored at the same set address.

One aspect of the present invention is a computer system comprising a main memory, a processor, a bus, a cache, and a protection bit. The computer system may be any system including, but not limited to, the systems shown in Figures 1, 6A-6C, 7, or 8. The bus connects the main memory and the processor. The cache is associated with the processor and has a plurality of cache lines. The protection bit is associated with each of the cache lines in each of the plurality of caches. Each protection bit protects a cache line from premature eviction during speculative execution. In one embodiment, the cache is a level one (L1) cache and in another embodiment, the cache is a level two (L2) cache. In one embodiment, the L1 cache is on the same chip die as the processor.

It is to be understood that the above description it is intended to be illustrative, and not restrictive. Many other embodiments will be apparent to those of skill in the art upon reviewing the above description. The scope of the invention should, therefore, be determined with reference to the appended claims, along with the full scope of equivalents to which such claims are entitled.